

Weakly-Supervised Spatial Context Networks

Zuxuan Wu¹, Larry S. Davis¹, Leonid Sigal²

¹University of Maryland, College Park, ² Disney Research

{zxwu, lsd}@umiacs.umd.edu, lsigal@disneyresearch.com

Abstract

We explore the power of spatial context as a self-supervisory signal for learning visual representations. In particular, we propose spatial context networks that learn to predict a representation of one image patch from another image patch, within the same image, conditioned on their real-valued relative spatial offset. Unlike auto-encoders, that aim to encode and reconstruct original image patches, our network aims to encode and reconstruct intermediate representations of the spatially offset patches. As such, the network learns a spatially conditioned contextual representation. By testing performance with various patch selection mechanisms we show that focusing on object-centric patches is important, and that using object proposal as a patch selection mechanism leads to the highest improvement in performance. Further, unlike auto-encoders, context encoders [21], or other forms of unsupervised feature learning, we illustrate that contextual supervision (with pre-trained model initialization) can improve on existing pre-trained model performance. We build our spatial context networks on top of standard VGG_19 and CNN_M architectures and, among other things, show that we can achieve improvements (with no additional explicit supervision) over the original ImageNet pre-trained VGG_19 and CNN_M models in object categorization and detection on VOC2007.

1. Introduction

Recent successful advances in object categorization, detection and segmentation have been fueled by high capacity deep learning models (e.g., CNNs) learned from massive labeled corpora of data (e.g., ImageNet [24], COCO [15]). However, the large-scale human supervision that makes these methods effective at the same time, limits their use; especially for fine-grained object-level tasks such as detection or segmentation, where annotation efforts become costly and unwieldily at scale. One popular solution is to use a pre-trained model (e.g., VGG_19 trained on ImageNet) for other, potentially unrelated, image tasks. Such pre-trained models produce effective and highly generic

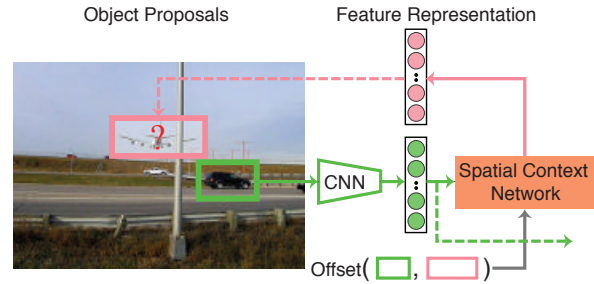


Figure 1. **Illustration of the proposed spatial context network.** A CNN used to compute feature representation of the green patch is fine-tuned to predict feature representation of the red patch using the proposed spatial context module, conditioned on their relative offset. Pairs of patches used to train the network are obtained from object proposal mechanisms. Once the network is trained, the green CNN can be used as a generic feature extractor for other tasks (dotted green line).

feature representations [4, 22]. However, it has also been shown that fine-tuning with task-specific labeled samples is often necessary [8].

Unsupervised learning is one way to potentially address some of these challenges. Unfortunately, despite significant research efforts unsupervised models such as auto-encoders [12, 29] and, more recently, context encoders [21] have not produced representations that can rival pre-trained models (let alone beat them). Among the biggest challenges is how to encourage a representation that captures semantic-level (e.g., object-level) information without having access to explicit annotations for object extent or class labels.

In the text domain, the idea of local spatial context within a sentence, proved to be an effective supervisory signal for learning distributed word vector representations (e.g., continuous bag-of-words (CBOW) [17] and skip-gram models [17]). The idea is conceptually simple; given a word tokenized corpus of text, to learn a representation for a target word that allows it to predict representations of contextual words around it; or vice versa, given contextual words to predict a representation of the target word. Generalizing this idea to images, while appealing, is also challenging as it

is not clear how to 1) *tokenize* the image (*i.e.*, what is an elementary entity between which context supervision should be applied) and 2) apply the notion of context effectively in a 2-dimensional real-valued domain.

Recent attempts to use spatial context as supervision in vision, resulted in models that used (regularly sampled) image patches as *tokens* and either learned a representation that is useful for classifying contextual relationships between them [3] or attempted to learn representations that fill in an image patch based on the larger surrounding pixels [21]. In both cases, the resulting feature representations fail to perform at the level of the pre-trained ImageNet models. This could be attributed to a number of reasons: 1) spatial context may indeed not be a good supervisory signal; 2) generic and neighboring image patches may not be an effective *tokenization* scheme; and/or 3) it may be difficult to train a model with a contextual loss from scratch.

Our motivation is similar to [3, 21]; however, we propose that image *tokenization* is important and should be done at the level of objects. By working with patches at object scale, our network can focus on more object-centric features and potentially ignore some of the texture and color details that are likely less important for semantic tasks. Further, instead of looking at immediate regions around the patch for context [21] and encoding the relationship between the contextual and target regions implicitly, we look at potentially non-overlapping patches with longer spatial contextual dependencies and explicitly condition the predicted representation on the relative spatial offset between the two regions. In addition, when training our network, we make use of a pre-trained model to extract intermediate representations. Since lower levels of CNNs have been shown to be task independent, this allows us to learn a better representation.

Specifically, we propose a novel architecture, spatial context network (SCN), which is built on top of existing CNN networks and is designed to predict a representation of one (object-like) image patch from another (object-like) image patch, conditioned on their relative spatial offset. As a result, the network learns a spatially conditioned contextual representation of image patches. In other words, given the same input patch and different spatial offsets it learns to predict different contextual representations (*e.g.*, given a patch depicting a side-view of a car and a horizontal offset, the network may output a patch representation of another car; however, the same input patch with a vertical offset may result in a patch representation of a plane). We also make use of ImageNet pre-trained model as both initialization and to define intermediate representation. Once SCN model is trained (on pairs of patches), we can use one of the two streams as a image representation that can be used for a variety of tasks, including object categorization or localization (*e.g.*, as part of Faster R-CNN [7]). This setting allows us to definitively answer the question of whether spa-

tial context can be an effective supervisory signal – it can, improving on the original ImageNet pre-trained models.

Contributions: Our main contribution is the spatial context network (SCN), which differs from other models in that it uses offset between two patches as a form of contextual supervision. Further, we explore a variety of tokenization schemes for mining training patch pairs, and show that an object proposal mechanism is the most effective. This observation validates the intuition that for semantic tasks, context is most useful at the object scale. Finally, we conduct extensive experiments to investigate the capacity of the proposed SCN for capturing context information in images, and demonstrate its ability to improve, in an unsupervised manner, on ImageNet pre-trained CNN models for both categorization (on VOC2007 and VOC2012) and detection (on VOC2007), where the bottom stream of the trained SCN is used as a generic feature extractor.

2. Related Work

Unsupervised Learning. Auto-encoders [11] are among the earliest works in unsupervised deep learning. Auto-encoders typically learn a representation by employing an encoder-decoder architecture; the encoder encodes the image (or patch) into a compact hidden state representation and the decoder reconstructs it back to a full image. Denoising auto-encoders [29] reconstruct images (or patches) subject to local corruptions. The most extreme variant of de-noising auto-encoders are the context encoders [21], which aim to reconstruct a large hole (patch) given its surrounding spatial context.

A number of papers proposed to learn representations by converting the generative auto-encoder-like objectives to discriminative classification counterparts, where CNNs have been shown to learn effectively. For example, [5] proposed an idea of surrogate classes that are formed by applying a variety of transformations to randomly sampled image patches. Classification into these surrogate classes is used as a supervisory signal to learn image representations. Alternatively, in [3], neighboring patches are used in Siamese-like networks to predict the relative *discrete* (*e.g.*, to the top-right, bottom-left, *etc.*) location of patches. Related, is also [34] that attempts to learn a similarity function across patches using various deep learning architectures, including center-surround (similar to [21]) and forms of Siamese networks. Goodfellow *et al.* [9] proposed Generative Adversarial Networks (GAN) that contain a generative model and discriminative model. Pathak *et al.* [21] built upon GANs to model context through inpainting missing patches.

Our model is related to auto-encoders [11], and particularly context encoders [21], however, it is conceptually somewhere between the discriminative and generative forms discussed above. We have encoder and decoder com-

ponents, but instead of decoding the hidden state all way to an image, our decoder decodes it to an intermediate discriminatively trained representation. Further, unlike previous methods, our decoder takes real-valued patch offsets as input, in addition to the representation of the patch itself.

Pre-trained Models. Pre-trained CNN models have been shown to generalize to a large number of different tasks [4, 22]. However, their transferability, as was noted in [33], is affected by specialization of higher layer neurons to the original task (often ImageNet categorization). By taking a network pre-trained on the ImageNet task and using its intermediate representation as target for our decoder, we make use of the knowledge distilled in the network [10] while attempting to improve it using spatial context. Works like [19] and [13] attempt to similarly re-use lower layers [19] of the pre-trained network and fine-tune, typically, fully-connected layers to specific tasks (e.g., object detection). However, such models assume some labeled data in the target domain, if not for classes of interest [19], then for related ones [13]. In our case, we assume no supervision of this form. Instead, we just assume that there exists a process that can generate category agnostic object-like proposal patches. Our work is similar to [37] that also attempts to improve the performance of pre-trained models. While they augment existing networks with reconstructive decoding pathways for image reconstruction, our model focuses on exploiting contextual cues in images.

Weakly-supervised and Self-supervised Learning. Recent years have witnessed a growing trend in weakly-supervised and self-supervised learning, which attempt to achieve similar performance to fully supervised models with limited use of annotated labels. A typical setting is to, for example, use image-level annotations to learn an object detection model [2, 20, 25, 27, 28, 31]. However, such models typically rely on latent variables and appearance regularities present within individual object class. In addition, researchers also utilized motion coherence (tracked patches [32] or ego-motion from sensors [1]) in videos as supervisory signals to train networks. Zhang *et al.* [36] proposed to generate a color version of a grayscale photo through a CNN model, which could further serve as an auxiliary task for feature learning. Noroozi *et al.* proposed to learn features by solving jigsaw puzzles [18]. Different from these works, we experiment with (category-independent) object proposals as a way to *tokenize* an image into more semantically meaningful parts. This can be thought of as (perhaps) a very weak form of supervision, but unlike any that we are aware has been used before.

Also related is [30], where the model for predicting future frame representation in video, given the current frame representation, is learned. The premise in [30] is conceptually similar to ours, but there are important differences. Our predictions are on spatial category-independent object

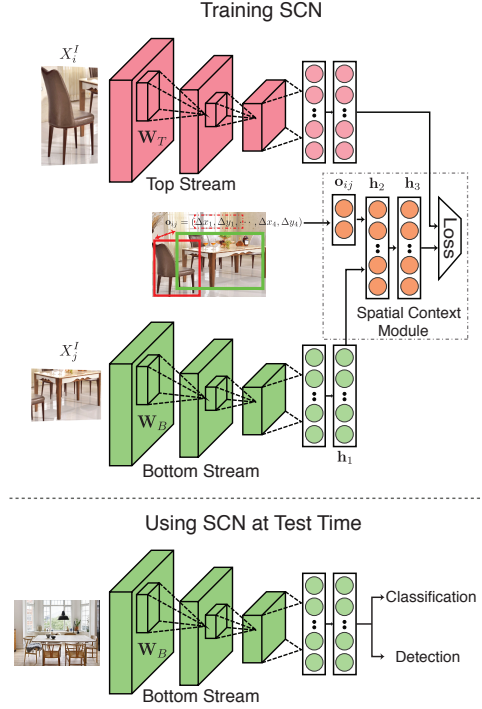


Figure 2. **Overview of the proposed spatial context network architecture.** See texts for complete description and discussion.

proposals (not frames offset in time [30]). Further, our neural network architecture is parametrized by the real-valued offset between pairs of proposals, where as temporal offset in [30] is not part of the model and is fixed to 1 second.

3. Spatial Context Networks

We now introduce the proposed spatial context network (see Figure 2 (top)), which consists of a top stream and a bottom stream operating on a pair of patches cropped from the same image. The goal is to utilize their spatial layout information as contextual clues for feature representation learning. Once the spatial context network is learned, the bottom stream can be used as a feature extractor (see Figure 2 (bottom)) for a variety of image recognition tasks, specifically, object categorization and detection.

More formally, given a patch \mathbf{X}_i^I extracted from the original image $I \in \mathbb{I}$, where \mathbb{I} is the training set, we denote the patch bounding box \mathbf{b}_i^I as an eight-tuple consisting of (x, y) positions of top-left, top-right, bottom-left and bottom-right corners. We can then denote the training samples for the network as 3-tuples $(\mathbf{X}_i^I, \mathbf{X}_j^I, \mathbf{o}_{ij}^I)$, where $\mathbf{o}_{ij}^I = \mathbf{b}_i^I - \mathbf{b}_j^I$ is the relative offset between two patches computed by subtracting locations of their respective four corners.

Top stream. The goal of the top stream is to provide a feature representation for patch \mathbf{X}_i^I that will be used as soft *tar-*

get for contextual prediction by the *learned* representation of the patch \mathbf{X}_j^I . This stream consists of an ImageNet pre-trained state-of-the-art CNN such as VGG_19, GoogleNet or ResNet (any pre-trained CNN model can be used). More specifically, the output of the top stream is the representation from the fully-connected layer (f_{c7}) obtained by propagating patch \mathbf{X}_j^I through the original pre-trained ImageNet model (here we remove the softmax layer). More formally, let $g(\mathbf{X}_j^I; \mathbf{W}_T)$ denote the non-linear function approximated by the CNN model and parameterized by weights \mathbf{W}_T . Note that one can also utilize representation of other layers; we use f_{c7} for simplicity and because of its superior performance in most high-level visual tasks [22].

Bottom stream. The bottom stream consists of an identical CNN model to the top stream which feeds into the proposed spatial context module. The spatial context module then accounts for spatial offset between the input pair of patches. The network first maps the input patch to a feature representation $\mathbf{h}_1 = g(\mathbf{X}_j^I; \mathbf{W}_B)$ and then the resulting \mathbf{h}_1 (f_{c7} representation) is used as input for the spatial context module. We initialize the bottom stream with the ImageNet pre-trained model as well, so initially, $\mathbf{W}_B = \mathbf{W}_T$.

Spatial Context Module. The role of the spatial context module is to take the feature representation of the patch \mathbf{X}_j^I produced by the bottom stream and, given the offset to patch \mathbf{X}_i^I , predict the representation of patch \mathbf{X}_i^I that would be produced by the top stream. The spatial context module is represented by a non-linear function $f([\mathbf{h}_1, \mathbf{o}_{ij}]; \mathbf{V})$, parameterized by weight matrix $\mathbf{V} = \{\mathbf{V}_1, \mathbf{V}_{loc}, \mathbf{V}_2\}$.

In particular, the spatial context module first takes the feature vector \mathbf{h}_1 (computed from patch \mathbf{X}_j^I) together with the offset vector \mathbf{o}_{ij} between \mathbf{X}_j^I and \mathbf{X}_i^I to derive an encoded representation:

$$\mathbf{h}_2 = \sigma(\mathbf{V}_1 \mathbf{h}_1 + \mathbf{V}_{loc} \mathbf{o}_{ij}), \quad (1)$$

where \mathbf{V}_1 denotes the weights for \mathbf{h}_1 ; \mathbf{V}_{loc} is the weight matrix for the input offset, and $\sigma(x) = 1/(1 + e^{-x})$. (Note that we absorb the bias term in the weight matrix for convenience). Finally, \mathbf{h}_2 is mapped to \mathbf{h}_3 with a linear transformation to reconstruct the f_{c7} feature vector computed by the top stream on the patch \mathbf{X}_i^I .

Loss Function. Given the output feature representations from the aforementioned two streams, we train the network by regressing the features from the bottom stream to those from the top stream. We use a squared loss function:

$$\min_{\mathbf{V}, \mathbf{W}_B} \sum_{I \in \mathbb{I}; i \neq j} \|g(\mathbf{X}_i^I; \mathbf{W}_T) - f([g(\mathbf{X}_j^I; \mathbf{W}_B), \mathbf{o}_{ij}]; \mathbf{V})\|^2. \quad (2)$$

The model is essentially an encoder-decoder framework with the bottom stream encoding the input image patch into a fixed representation and spatial context module decoding it to representation of another, spatially offset, patch.

The intuition comes from the skip-gram model [16] that attempts to predict the context given a word, which has been demonstrated to be effective for a number of NLP tasks. Since objects often co-occur in images in particular relative locations, it makes intuitive sense to explore such relations as contextual supervision.

The network can be easily trained using back-propagation with stochastic gradient descent. Note that for the top stream, rather than predicting raw pixels in images, we utilize the features extracted from off-the-shelf CNN architecture as *ground truth*, to which the features constructed by the bottom stream regress. This is because the pre-trained CNN model contains valuable semantic information (e.g., referred to as dark knowledge [10]) to differentiate objects and the extracted off-the-shelf features have achieved great success on various tasks [35, 38].

One alternative to formulating the problem as a regression task, would be to turn it into a classification problem by appending a softmax layer on top of the two streams and predicting whether a pair of features is likely given the spatial offset. However, this would require a large amount of negative samples (e.g., a *car* is not likely to be in a *lake*), making training difficult. Further, our regression loss also builds on intuitions explored in [10], where it is shown that soft real-valued targets are often better than discrete labels.

Implementation Details. We adopt two off-the-shelf CNN architectures, CNN_M and VGG_19 [26], to train the spatial context network. CNN_M is an AlexNet [14] style CNN with five convolutional layers topped by three fully-connected layers (the dimension for f_{c6} and f_{c7} is 2,048), but contains more convolutional filters. VGG_19 network consists of 16 convolutional layers followed by three fully-connected layers, possessing stronger discriminative power.

The pipeline was implemented in Torch and we apply mini-batch stochastic gradient descent in training with the batch size of 64. The weights for the spatial context module are initialized randomly. We fine-tune the fully-connected layers in the bottom stream CNN model with convolutional layers fixed, unless otherwise specified. The input patches are resized to 224×224 . We set the initial learning rate to $1e^{-3}$, which is decreased to $1e^{-4}$ after 100 epochs; we fix weight decay to $5e^{-4}$ and the maximum number of epochs to 200. We will discuss patch selection in Experiments.

3.1. Using SCN for Classification and Detection

Once the SCN is trained, we use \mathbf{h}_1 from the bottom stream as a feature representation for other tasks (Figure 2 (bottom)). As we will show, these feature representations are better than those obtained from the original ImageNet pre-trained model for object detection and classification.

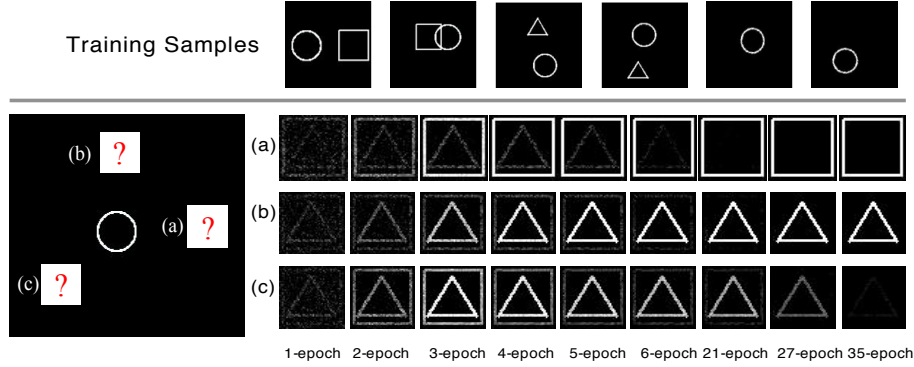


Figure 3. **Experiments with synthetic dataset.** Training samples are shown in top row. Bottom rows show predicted patches for the labeled regions on the left, after 1–35 epochs of training. Predicted patches are obtained by treating the circle in the middle and an appropriate spatial offset to (a), (b), or (c) as input to an SCN and visualizing the output h_3 layer.

4. Experiments

We first validate the ability of the proposed SCN to learn context information on a synthetic dataset and with the real images from VOC2012. We then evaluate the effectiveness of features extracted from the spatial context framework in classification and detection tasks, as compared with original pre-trained ImageNet features, and competing state-of-the-art feature learning methods.

4.1. Synthetic Dataset Experiments

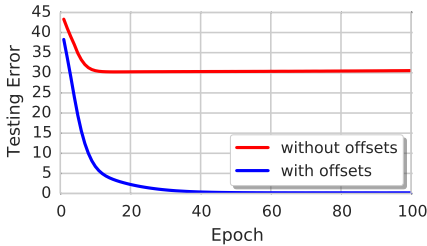


Figure 4. **Testing error on the synthetic dataset.** Illustrated is the testing error with and without offset vector in the input.

We construct a synthetic dataset containing *circles*, *squares* and *triangles* to verify whether the proposed spatial context framework is able to learn correlations in spatial layout patterns of these objects. More specifically, we create 300 (circle, square) pairs where circles are always horizontally offset (see Figure 3 (top)) from the squares (vertical difference is within 30 pixels); and 300 (circle, triangle) pairs where circles are vertically offset from the triangles (horizontal difference is within 30 pixels); as well as 200 (circles, black image) pairs where the offset vector is randomly sampled. We randomly split the dataset into 600 training and 200 testing pairs. We assume perfect proposals and crop patches tightly around the objects (circles, squares and triangles). Here, we adopt the CNN_M model only.

The testing error loss (mean squared error) on this dataset is visualized in Figure 4. As we can see from the fig-

ure, the testing error of the spatial context network steadily decreases for the first 20 epochs and nearly reaches zero after 25 epochs. To investigate the role offset vectors play in the learning process, we remove the offset vector from the input and retrain the network. The loss of this network stabilizes to 30 after 10 epochs; this is significantly higher than the error of the spatial context network. Figure 4 confirms that the proposed spatial context network can make effective use of the spatial context information between objects.

To gain further insights into the learning process, we replace the *target* feature representation of the top stream with raw ground truth image patches. After each epoch, given an input bottom stream object patch (depicting circle) and an offset vector from the testing set, we adopt the output of the last layer h_3 in the SCN to reconstruct images for the top stream. The results are visualized in Figure 3 (bottom).

When circles are combined with either horizontal or vertical offsets, the network is able to reconstruct square and triangle patches (respectively) after about five epochs. For the first few epochs, both triangles and squares co-occur in the constructed images, but clear square and triangle patterns emerge as the training proceeds. It took longer for the network to learn that conditioned on an off-axis offset vector and a circle patch it should produce empty (black) patch image. This experiment validates that our spatial context network is able to learn correct spatially varying contextual representation based on (identical) input patch (circle) and varying offsets. Without providing location offset information, the network overfits and simply generates a patch containing overlapping triangles and squares (which explains poor convergence in Figure 4).

Imagining a circle is a car, a square a tree and the triangle (which is above circle) to be sky, this synthetic dataset provides a simplified version of spatial context information in real-world scenarios. The experiments indicate that the varying spatial contextual information among multiple objects can be learned by the SCN.

4.2. Modeling Context in Real Images

We now discuss context modeling in real images and validate the capability of the network to capture such real-world contextual clues. To this end, we adopt PASCAL VOC 2012 [6] dataset, which consists of a training set with 5,717 images and a validation set with 5,823 images, totaling 20 object categories (denoted by VOC2012-Img). We first crop objects from the original images on both subsets using the provided annotations of bounding boxes, which leads to 15,774 objects for training and 15,787 objects for testing (denoted by VOC2012-Obj¹). Objects from the same image are further paired and are used as inputs for the spatial context network (SCN) together with their offset vector. In total, we obtain 34,378 training and 34,722 testing paired samples (VOC2012-Pairs).

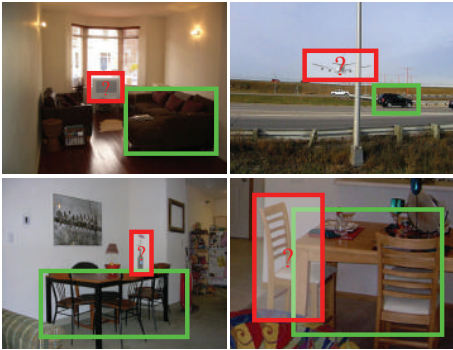


Figure 5. **SCN contextual classification.** Features of the top stream (red boxes) are predicted using patches from bottom stream (green boxes) and offset vector as inputs to the trained SCN. A classifier is then trained to predict the label of the red patch based on the *predicted* features from the training set. Performance on testing set is 56.3% (Table 1).

We first train the spatial context network using paired images. Given the trained network, we compute the outputs of the last layer from the spatial context module (*i.e.*, \mathbf{h}_3) as the synthesized/predicted feature representations for a single patch in the top stream (on both training and test set). Then we train a linear classifier with the extracted features using all training patches in the top stream (See Figure 5 for illustration). To establish a baseline, for all patches in the top stream, we compute the raw f_{c7} features from the original VGG-19 network and similarly train a linear SVM classifier. The results are summarized in Table 1.

It is surprising to see that the predicted features achieve a 56.3% accuracy in object classification given the fact that these features are *predicted* from nearby objects within the same image (from the bottom stream) using the trained spatial context network (SCN). In other words we are able to recognize objects at 56.3% accuracy without ever seeing

¹The difference between VOC2012-Obj and VOC2012-Img is that in the former the objects are cropped, where as in the latter they are not.

features	VOC2012-Pairs (%)
VGG-19 f_{c7}	78.3
SCN predicted (\mathbf{h}_3) features	56.3
VGG-19 f_{c7} + SCN predicted	79.5

Table 1. **Performance comparisons of classification.** Different feature representations for the *top* patch classification are compared. SCN predicted features are obtained by regressing *top* stream features from the contextual *bottom* stream patch.

the real image features contained in the corresponding image patches; the recognition is done purely based on the *contextual* predictions of those features from other patches (note that 92.6% of patches do not or minimally overlap (< 0.2 IoU)). This indicates *very strong* contextual information that our network was able to learn.

To eliminate the possibility that accuracy comes from images containing multiple instances of the same object, we analyzed the dataset and found only 45% of training and 42% of testing image patch pairs correspond to the same objects. Further, using pairs that do not contain same objects produces an accuracy of 52.8%, and 63.2% with pairs only from the same objects.

To investigate whether the synthesized features \mathbf{h}_3 contain contextual information that might be complementary to the original f_{c7} features, we perform feature fusion by concatenating the two representations into a 8,192-dimensional vector and training a linear SVM for classification. We observe 1.2% performance gain compared with raw VGG f_{c7} features, which again confirms that context is beneficial.

4.3. Feature Learning with SCN for Classification

In the last two sections, to verify the effectiveness of spatial contextual learning, we assumed knowledge of object bounding boxes (but, importantly, not their categorical identity); in other words, we assumed existence of a perfect object proposal mechanism; this is clearly unrealistic. In this section, we explore the importance/significance of the quality of the object proposal mechanism on the performance of features learned using SCN. We do so in the context of classification, where once SCN is trained, we use SVM on top of generic SCN features (see Figure 2 (bottom)).

We use ground truth bounding boxes, provided by the dataset, as a baseline (*SCN-BBox*). In addition, we test the following object proposal methods:

- **Random Patches** (*SCN-Random*): We randomly crop 5 patches of size of 64×64 in each image (consistent with [21]) to generate 10 patch pairs per image. In total, we collect 28K cropped patches and 57K pairs.²
- **Edge Box** [39] (*SCN-EdgeBox*): EdgeBox is a generic method to generate object bounding box proposals

²Note that in the pairing process one could simply swap the inputs of the top and bottom stream to double the number of pairs for the network, however, empirically, we found it not to be helpful.

	features- f_{c7}	VOC2012-Obj	VOC2012-Img
CNN_M	Original	75.3	68.5
	SCN-BBox	78.7	70.8
	SCN-YOLO	79.2	70.7
	SCN-EdgeBox	79.9	72.8
	SCN-Random	78.8	70.0
VGG_19	Original	81.4	78.1
	SCN-BBox	82.6	78.8
	SCN-YOLO	83.0	79.0
	SCN-EdgeBox	83.6	79.5
	SCN-Random	83.2	79.2

Table 2. **Performance with various object proposals.** Comparison of classification with features obtained using SCN trained with different patch selection mechanisms is illustrated on VOC2012-Obj and VOC2012-Img, using two CNN architectures.

based on edge responses. We filter out the bounding boxes with confidence lower than 0.1 and those with irregular aspect ratio, leading to 43K *object* patches and 160K pairs for training.

- **YOLO [23] (SCN-YOLO):** YOLO is a recently introduced end-to-end framework trained on VOC for object detection. We use YOLO as an object proposal mechanism, by taking patches from detection regions but ignoring the detected labels. We collect 13K objects forming 17K image patch pairs.

We expect the quality of object proposal methods (from least object-like to most object-like) on VOC to roughly follow the following pattern:

$$Random < EdgeBox < YOLO < \text{ground-truth } BBox.$$

Given a trained SCN model, we utilize the bottom stream (see Fig. 2 (bottom)) to test generalization of the learned feature representations, by performing classification with linear SVMs on VOC2012-Obj and VOC2012-Img (see footnote 1 for explanation) with the outputs from the first hidden layer (h_1 , *i.e.*, fine-tuned version of f_{c7}) in the bottom stream of SCN. The results are measured in mAP. We compare the different patch selection mechanisms discussed above and also to the original ImageNet pre-trained models. The results are summarized in Table 2. We observe that SCN-BBox and SCN-YOLO achieve better results compared with the original f_{c7} features. It is also surprising to see that SCN-EdgeBox obtains the best performance, even higher than models trained with ground-truth bounding boxes. SCN-EdgeBox is 4.6 and 4.3 percentage points better than the original f_{c7} features on VOC2012-Obj and VOC2012-Img respectively.

We believe that better performance of the SCN-EdgeBox stems from EdgeBox’s ability to select object-like regions that go beyond the 20 object classes labeled in ground truth and detected by YOLO. We also note that while Random

	VOC2012-Obj
VGG_19 f_{c7}	81.4
SCN-EdgeBox (f_{c6}, f_{c7})	83.6
SCN-EdgeBox ($f_{c6}, f_{c7}, conv_5$)	84.3
SCN-EdgeBox (all layers)	82.5

Table 3. **Exploring SCN learning strategies.** Classification performance based on features obtained using different fine-tuning strategies. See text for more details.

patch sampling also improves the performance, with respect to the original ImageNet pre-trained network, it is doing so by a much smaller margin than EdgeBox patch sampling.

The original f_{c7} features are trained using labels from ImageNet; our spatial context network is appealing in that it learns a better feature representation by exploiting contextual cues without any additional explicit supervision. Figure 6 compares the per-class performance of SCN-EdgeBox and the original f_{c7} features on VOC2012-Img, where we can see that SCN-EdgeBox features outperform the original f_{c7} features for all classes. It is also interesting to see that, for small objects, such as “bottle” and “potted plant”, the performance gain of SCN-EdgeBox is more significant.

Fine-tuning Convolutional Layers. In addition to only fine-tuning the fully-connected layers of the bottom stream CNN model, we also explore whether joint training with VGG_19 network could further improve the performance of the extracted features. More specifically, for the top stream we fix the weights since computing features dynamically poses challenges for network convergence. Further, this avoids *trivial* solutions of both streams learning, for example, to predict zero features for all patches. In addition, this makes use of transferability of lower levels of pre-trained CNN models as targets for the bottom stream decoding. The results are summarized in Table 3. By back-propagating the error through deeper layers we observe a significant performance gain (2.9 percentage points) over the original features of VGG_19 network, which confirms the fact that SCN is effective and VGG layers could be fine-tuned jointly for specific tasks in order to gain better performance using our formulation. When fine-tuning all layers in the network, the performance of SCN degrades slightly to 82.5%.

4.4. Feature Learning with SCN for Detection

We also explore the applicability of SCN features for object detection tasks to verify generic feature effectiveness. To make fair comparisons with prior work, we adopt the experimental setting of [21] and fine-tune the SCN-EdgeBox model (based on CNN_M architecture) on Pascal VOC2007, which is then applied in the *Fast R-CNN* [7] framework. More precisely, we replace the ImageNet pre-trained CNN_M model with the fine-tuned bottom stream in SCN (See Figure 2 (bottom)). The weights for final classi-

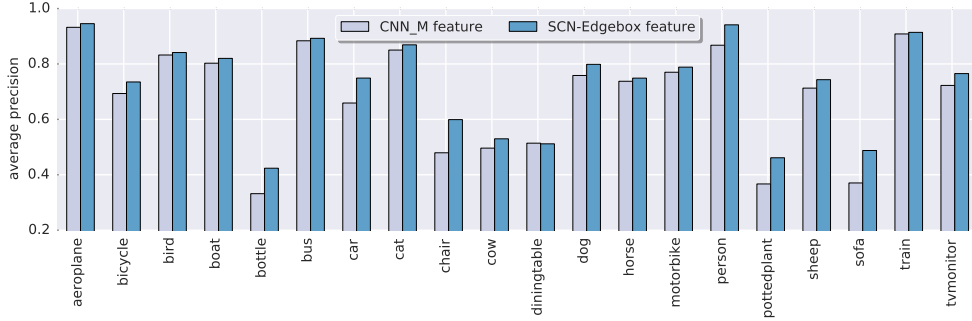


Figure 6. **Classification per class performance.** Reported is average precision obtained using original CNN_M features and SCN-EdgeBox features on VOC2012-Img.

	Initialization	Supervision	Pretraining time	Classification	Detection
Random Gaussian	random	N/A	< 1 minute	53.3	43.4
Wang <i>et al.</i> [32]	random	motion	1 week	58.4	44.0
Doersch <i>et al.</i> [3]	random	context	4 weeks	55.3	46.6
*Doersch <i>et al.</i> [3]	1000 class labels	context	–	65.4	50.4
Pathak <i>et al.</i> [21]	random	context inpainting	14 hours	56.5	44.5
Zhang <i>et al.</i> [36]	random	color	–	65.6	46.9
ImageNet [21]	random	1000 class labels	3 days	78.2	56.8
*ImageNet	random	1000 class labels	3 days	76.9	58.7
SCN-EdgeBox	1000 class labels	context	10 hours	79.0	59.4

Table 4. **Quantitative comparison for classification and detection on VOC 2007.** Classification and Fast-RCNN detection results are on the PASCAL VOC 2007 test set. The baselines labeled with * are based on our experiments, rest taken from original papers.

fication and bounding box regression layers are initialized from scratch. We then follow the training and testing protocol defined in [7] and report detector performance in mAP.

The results and comparisons with existing state-of-the-art methods are summarized in Table 4. SCN-EdgeBox model improves on the original ImageNet pre-trained model by 0.7 percentage points. Further, compared with alternative unsupervised learning methods, our approach achieves significantly better performance. We also significantly outperform other feature training methods on classification (including our fine-tuned ImageNet model) and Doersch *et al.* [3] model initialized with ImageNet.

Figure 7 visualizes some sample images where SCN-EdgeBox outperforms the pre-trained ImageNet model. Our model is better at detecting relatively small objects (*e.g.*, airplane in the first row and chair in the second row).

5. Conclusion

In this paper, we present a novel spatial context network built on top of existing CNN architectures. SCN network exploits implicit contextual layout cues in images as a supervisory signal. More specifically, the network is trained to predict the intermediate representation of one (object-like) image patch from another (object-like) image patch, within the same image, conditioned on their relative spatial offset. Consequently, the network learns a spatially conditioned

contextual representation of image patches. Extensive experiments are conducted to validate the effectiveness of the proposed spatial context network in modeling context information in images. We show that the proposed spatial context network can achieve improvements (with no additional explicit supervision) over the original ImageNet pre-trained models in object categorization on VOC2007 / VOC2012 and detection on VOC2007.

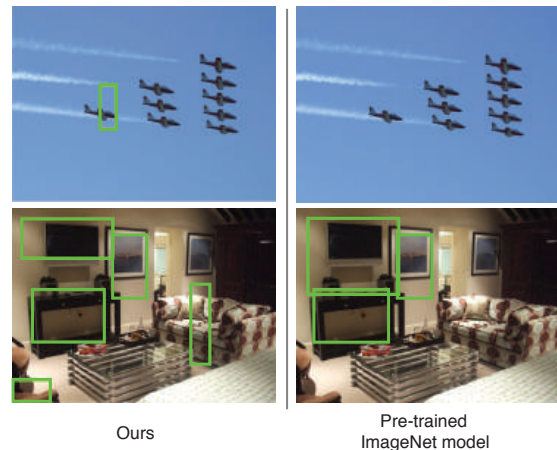


Figure 7. **Sample detection results.** Illustrated are results obtained using SCN-EdgeBox model and the original pre-trained ImageNet model, respectively, on VOC2007.

References

- [1] P. Agrawal, J. Carreira, and J. Malik. Learning to see by moving. In *ICCV*, 2015. 3
- [2] R. Cinbis, J. Verbeek, and C. Schmid. Multi-fold mil training for weakly supervised object localization. In *CVPR*, 2014. 3
- [3] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015. 2, 8
- [4] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *ICML*, 2014. 1, 3
- [5] A. Dosovitskiy, J. T. Springenberg, M. Riedmiller, and T. Brox. Discriminative unsupervised feature learning with convolutional neural networks. *TPAMI*, 2015. 2
- [6] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010. 6
- [7] R. Girshick. Fast r-cnn. In *ICCV*, 2015. 2, 7, 8
- [8] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 1
- [9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. 2
- [10] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *CoRR*, 2015. 3, 4
- [11] G. E. Hinton. Learning multiple layers of representation. *Trends in cognitive sciences*, 2007. 2
- [12] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 2006. 1
- [13] J. Hoffman, S. Guadarrama, E. Tzeng, R. Hu, J. Donahue, R. Girshick, T. Darrell, and K. Saenko. Lsda: Large scale detection through adaptation. In *NIPS*, 2014. 3
- [14] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 4
- [15] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár. Microsoft coco: Common objects in context. In *ECCV*, 2014. 1
- [16] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013. 4
- [17] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013. 1
- [18] M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016. 3
- [19] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *CVPR*, 2014. 3
- [20] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Is object localization for free? – weakly-supervised learning with convolutional neural networks. In *CVPR*, 2015. 3
- [21] D. Pathak, P. Krähenbühl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 1, 2, 6, 7, 8
- [22] A. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *CVPR workshop of DeepVision*, 2014. 1, 3, 4
- [23] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 7
- [24] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 1
- [25] Z. Shi, T. M. Hospedales, and T. Xiang. Bayesian joint topic modelling for weakly supervised object localisation. In *ICCV*, 2013. 3
- [26] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 4
- [27] H. O. Song, R. Girshick, S. Jegelka, J. Mairal, Z. Harchaoui, and T. Darrell. On learning to localize objects with minimal supervision. In *ICML*, 2014. 3
- [28] H. O. Song, Y. J. Lee, S. Jegelka, and T. Darrell. Weakly supervised discovery of visual pattern configurations. In *NIPS*, 2014. 3
- [29] P. Vincent, H. Larochelle, Y. Bengio, and P. A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *ICML*, 2008. 1, 2
- [30] C. Vondrick, H. Pirsiavash, and A. Torralba. Anticipating visual representations with unlabeled videos. In *CVPR*, 2016. 3
- [31] C. Wang, W. Ren, K. Huang, and T. Tan. Weakly supervised object localization with latent category learning. In *ECCV*, 2014. 3
- [32] X. Wang and A. Gupta. Unsupervised learning of visual representations using videos. In *ICCV*, 2015. 3, 8
- [33] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *NIPS*, 2014. 3
- [34] S. Zagoruyko and N. Komodakis. Learning to compare image patches via convolutional neural networks. In *CVPR*, 2015. 2
- [35] S. Zha, F. Luisier, W. Andrews, N. Srivastava, and R. Salakhutdinov. Exploiting image-trained cnn architectures for unconstrained video classification. In *BMVC*, 2015. 4
- [36] R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. In *ECCV*, 2016. 3, 8
- [37] Y. Zhang, K. Lee, and H. Lee. Augmenting supervised neural networks with unsupervised objectives for large-scale image classification. In *ICML*, 2016. 3
- [38] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *NIPS*, 2014. 4
- [39] C. L. Zitnick and P. Dollár. Edge boxes: Locating object proposals from edges. In *ECCV*, 2014. 7